

Machine Learning for Data Science (CS4786)  
Lecture 8

Kernel PCA  
&  
Isomap

# KERNEL TRICK

- We have have nice methods for linear dimensionality reduction
- Can we use this beyond the linear realm?

# KERNEL TRICK

- Lift to higher dimensions (introduces non-linearity)
  - Using feature map  $\phi(x)$
- Perform linear dimensionality reduction in this high dimensional space

# Key Idea:

If an algorithm only depends on inner products, we can simply replace inner product in  $x$  space by inner product in  $\phi(x)$  space

Can we write PCA so it only depends on inner products?

# LETS REWRITE PCA

# LETS REWRITE PCA

**Lets start with the assumption that Data is centered! (i.e. Sum of  $x_t$ 's is 0)**

# LETS REWRITE PCA

Lets start with the assumption that Data is centered! (i.e. Sum of  $\mathbf{x}_t$ 's is 0)

- $k^{\text{th}}$  column of  $W$  is eigenvector of covariance matrix  
That is,  $\lambda_k W_k = \Sigma W_k$ . Rewriting, for centered  $X$

$$\lambda_k W_k = \frac{1}{n} \left( \sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^{\top} \right) W_k = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t^{\top} W_k) \mathbf{x}_t$$



# LETS REWRITE PCA

Lets start with the assumption that Data is centered! (i.e. Sum of  $\mathbf{x}_t$ 's is 0)

- $k^{\text{th}}$  column of  $W$  is eigenvector of covariance matrix  
That is,  $\lambda_k W_k = \Sigma W_k$ . Rewriting, for centered  $X$

$$\lambda_k W_k = \frac{1}{n} \left( \sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^{\top} \right) W_k = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t^{\top} W_k) \mathbf{x}_t$$

But  $\mathbf{x}_t^{\top} W_k = \mathbf{y}_t[k]$

# LETS REWRITE PCA

Lets start with the assumption that Data is centered! (i.e. Sum of  $\mathbf{x}_t$ 's is 0)

- $k^{\text{th}}$  column of  $W$  is eigenvector of covariance matrix  
That is,  $\lambda_k W_k = \Sigma W_k$ . Rewriting, for centered  $X$

$$\lambda_k W_k = \frac{1}{n} \left( \sum_{t=1}^n \mathbf{x}_t \mathbf{x}_t^{\top} \right) W_k = \frac{1}{n} \sum_{t=1}^n (\mathbf{x}_t^{\top} W_k) \mathbf{x}_t$$

But  $\mathbf{x}_t^{\top} W_k = \mathbf{y}_t[k]$

$$\lambda_k W_k = \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t$$

# LETS REWRITE PCA

$$\mathbf{y}_s[k] = W_k^\top \mathbf{x}_s$$

# LETS REWRITE PCA

$$\begin{aligned}\mathbf{y}_s[k] &= W_k^\top \mathbf{x}_s \\ &= \frac{1}{\lambda_k} \left( \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \mathbf{x}_s\end{aligned}$$

# LETS REWRITE PCA

$$\begin{aligned}\mathbf{y}_s[k] &= W_k^\top \mathbf{x}_s \\ &= \frac{1}{\lambda_k} \left( \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \mathbf{x}_s \\ &= \frac{1}{n\lambda_k} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t^\top \mathbf{x}_s\end{aligned}$$

# LETS REWRITE PCA

$$\begin{aligned}\mathbf{y}_s[k] &= W_k^\top \mathbf{x}_s \\ &= \frac{1}{\lambda_k} \left( \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \mathbf{x}_s \\ &= \frac{1}{n\lambda_k} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t^\top \mathbf{x}_s \\ &= \frac{1}{n\lambda_k} \sum_{t=1}^n \mathbf{y}_t[k] \tilde{K}_{s,t}\end{aligned}$$

# LETS REWRITE PCA

$$\begin{aligned}\mathbf{y}_s[k] &= W_k^\top \mathbf{x}_s \\ &= \frac{1}{\lambda_k} \left( \frac{1}{n} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \mathbf{x}_s \\ &= \frac{1}{n\lambda_k} \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t^\top \mathbf{x}_s \\ &= \frac{1}{n\lambda_k} \sum_{t=1}^n \mathbf{y}_t[k] \tilde{K}_{s,t}\end{aligned}$$

Where  $\tilde{K}_{s,t} = \mathbf{x}_t^\top \mathbf{x}_s$  is the kernel matrix for centered data

# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$



# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$

Also we have,  $1 = \|W_k\|^2 = \frac{1}{\lambda_k^2 n^2} \left( \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \left( \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s \right)$

# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$

Also we have,  $1 = \|W_k\|^2 = \frac{1}{\lambda_k^2 n^2} \left( \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \left( \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s \right)$

$$= \frac{1}{\lambda_k^2 n^2} \sum_{t=1}^n \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s^\top \mathbf{x}_t \mathbf{y}_t[k]$$

# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$

Also we have,  $1 = \|W_k\|^2 = \frac{1}{\lambda_k^2 n^2} \left( \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \left( \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s \right)$

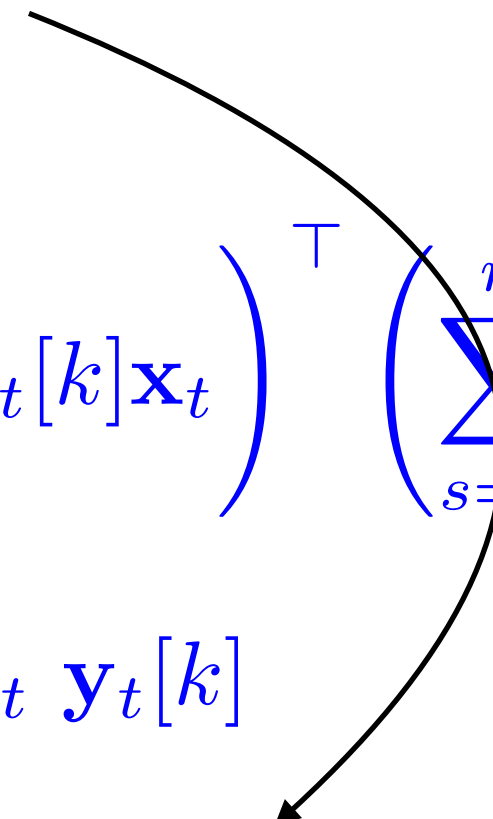
$$= \frac{1}{\lambda_k^2 n^2} \sum_{t=1}^n \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s^\top \mathbf{x}_t \mathbf{y}_t[k]$$
$$= \frac{1}{\lambda_k^2 n^2} \mathbf{y}[k] \tilde{K} \mathbf{y}[k]^\top$$

# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$

Also we have,  $1 = \|W_k\|^2 = \frac{1}{\lambda_k^2 n^2} \left( \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \left( \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s \right)$

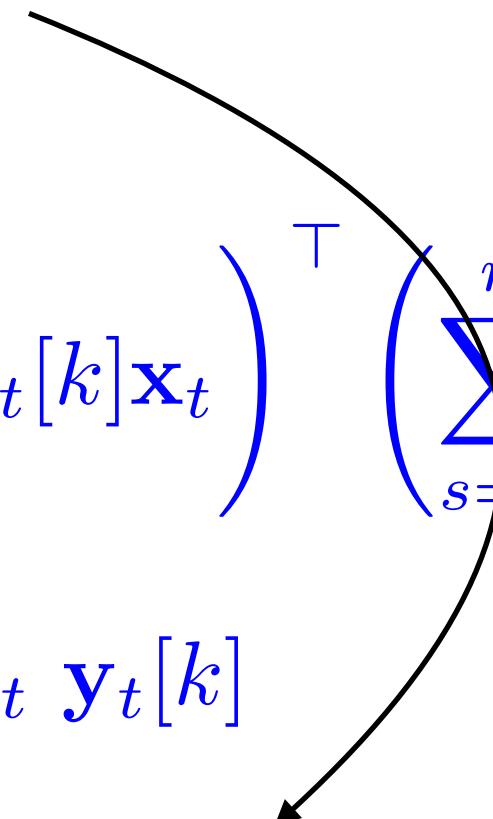
$$= \frac{1}{\lambda_k^2 n^2} \sum_{t=1}^n \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s^\top \mathbf{x}_t \mathbf{y}_t[k]$$
$$= \frac{1}{\lambda_k^2 n^2} \mathbf{y}[k] \tilde{K} \mathbf{y}[k]^\top = \frac{1}{n\lambda_k} \|\mathbf{y}[k]\|^2$$


# LETS REWRITE PCA

- Hence, the  $k$ 'th column on  $Y$  matrix is such that

$$\mathbf{y}[k] = \frac{1}{n\lambda_k} \mathbf{y}[k] \tilde{K}$$

Also we have,  $1 = \|W_k\|^2 = \frac{1}{\lambda_k^2 n^2} \left( \sum_{t=1}^n \mathbf{y}_t[k] \mathbf{x}_t \right)^\top \left( \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s \right)$

$$= \frac{1}{\lambda_k^2 n^2} \sum_{t=1}^n \sum_{s=1}^n \mathbf{y}_s[k] \mathbf{x}_s^\top \mathbf{x}_t \mathbf{y}_t[k]$$
$$= \frac{1}{\lambda_k^2 n^2} \mathbf{y}[k] \tilde{K} \mathbf{y}[k]^\top = \frac{1}{n\lambda_k} \|\mathbf{y}[k]\|^2$$


Hence  $P_k = \mathbf{y}[k] / \sqrt{n\lambda_k}$  is an eigenvector of  $\tilde{K}$  with eigen value  $\gamma_k = n\lambda_k$

# REWRITING PCA

- We assumed centered data, what if its not,

# REWRITING PCA

- We assumed centered data, what if its not,

$$\tilde{K}_{s,t} = \left( \mathbf{x}_t - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \left( \mathbf{x}_s - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)$$

# REWRITING PCA

- We assumed centered data, what if its not,

$$\begin{aligned}\tilde{K}_{s,t} &= \left( \mathbf{x}_t - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \left( \mathbf{x}_s - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right) \\ &= \mathbf{x}_t^\top \mathbf{x}_s - \left( \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \mathbf{x}_s - \left( \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \mathbf{x}_t \\ &\quad + \frac{1}{n^2} \left( \sum_{u=1}^n \mathbf{x}_u \right)^\top \left( \sum_{v=1}^n \mathbf{x}_v \right)\end{aligned}$$



# REWRITING PCA

- We assumed centered data, what if its not,

$$\begin{aligned}\tilde{K}_{s,t} &= \left( \mathbf{x}_t - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \left( \mathbf{x}_s - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right) \\ &= \mathbf{x}_t^\top \mathbf{x}_s - \left( \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \mathbf{x}_s - \left( \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u \right)^\top \mathbf{x}_t \\ &\quad + \frac{1}{n^2} \left( \sum_{u=1}^n \mathbf{x}_u \right)^\top \left( \sum_{v=1}^n \mathbf{x}_v \right) \\ &= \mathbf{x}_t^\top \mathbf{x}_s - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u^\top \mathbf{x}_s - \frac{1}{n} \sum_{u=1}^n \mathbf{x}_u^\top \mathbf{x}_t + \frac{1}{n^2} \sum_{u=1}^n \sum_{v=1}^n \mathbf{x}_u^\top \mathbf{x}_v\end{aligned}$$

# REWRITING PCA

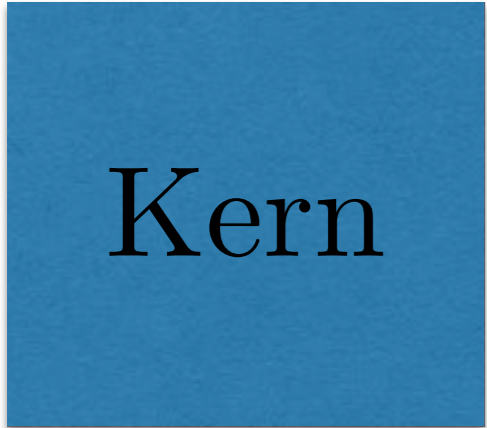
- Equivalently, if **Kern** is the matrix ( $\text{Kern}_{t,s} = x_t^\top x_s$ ),

$$\tilde{K} = \text{Kern} - \frac{(\mathbf{1}_{n \times n} \times \text{Kern})}{n} - \frac{(\text{Kern} \times \mathbf{1}_{n \times n})}{n} + \frac{(\mathbf{1}_{n \times n} \times \text{Kern} \times \mathbf{1}_{n \times n})}{n^2}$$

# KERNEL PCA

# KERNEL PCA

1.



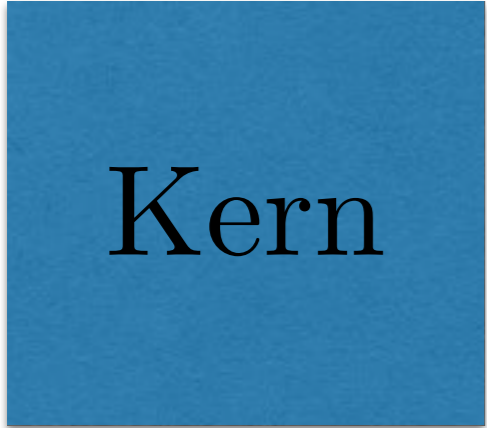
n

n

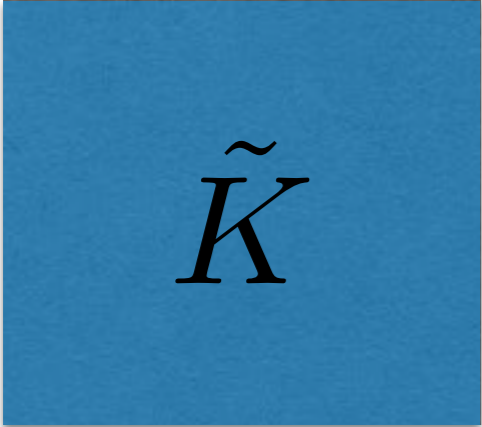
$$= \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ k(x_{n-1}, x_1) & k(x_{n-1}, x_2) & \dots & k(x_{n-1}, x_n) \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix}$$

# KERNEL PCA

1.


$$= \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \dots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \dots & k(x_2, x_n) \\ \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot \\ k(x_{n-1}, x_1) & k(x_{n-1}, x_2) & \dots & k(x_{n-1}, x_n) \\ k(x_n, x_1) & k(x_n, x_2) & \dots & k(x_n, x_n) \end{bmatrix}$$

2.


$$= \text{Kern} - \frac{1}{n} (\mathbf{1} \text{ Kern} + \text{Kern} \mathbf{1}) + \frac{1}{n^2} \mathbf{1} \text{ Kern} \mathbf{1}$$

# KERNEL PCA

# KERNEL PCA

$$3. \left[ \begin{array}{c} n \\ \mathbf{P} \\ K \end{array} , \gamma \right] = \text{eigs} \left( \begin{array}{c} \tilde{\mathbf{K}} \\ K \end{array} \right)$$

# KERNEL PCA

$$3. \begin{bmatrix} n \\ \mathbf{P} \\ K \end{bmatrix}, \gamma = \text{eigs} \left( \begin{bmatrix} \tilde{K} \\ K \end{bmatrix} \right)$$

$$4. \begin{bmatrix} n \\ \mathbf{Y} \\ K \end{bmatrix} = \begin{bmatrix} \vdots & \vdots \\ P_1 \sqrt{\gamma_1} & P_K \sqrt{\gamma_K} \\ \vdots & \vdots \end{bmatrix}$$

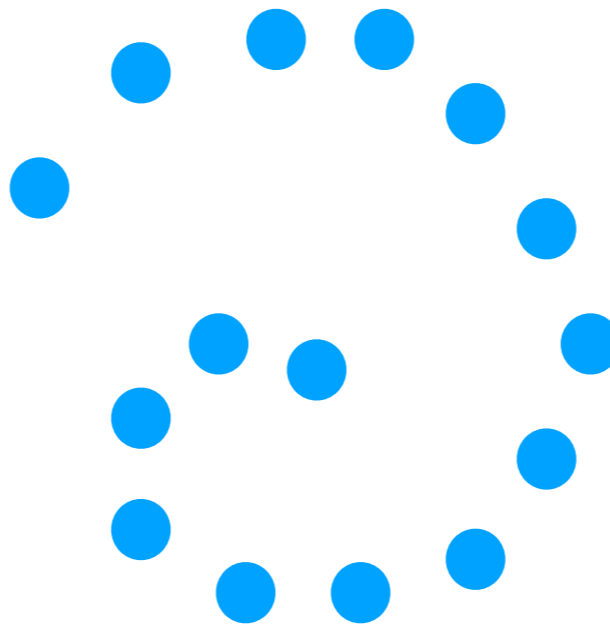


# Kernel Methods: A note

- We can “kernelize” CCA and many other linear dimensionality reduction methods.
- For any linear method, solution lies within linear span of data
- For typical linear methods  $y$ 's can be computed only based on inner products.

# MANIFOLD BASED DIMENSIONALITY REDUCTION

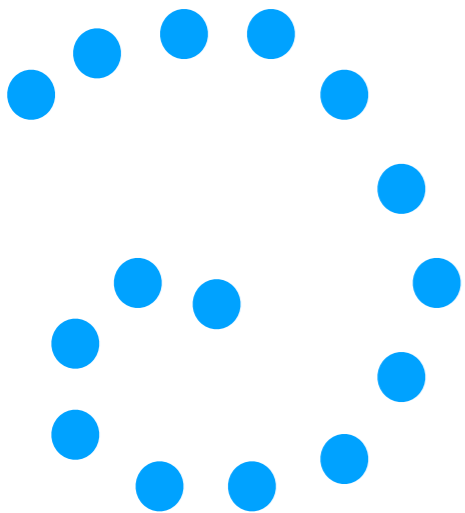
- Key Assumption: Points live on a low dimensional manifold
- Manifold: subspace that looks locally Euclidean
- Given data, can we uncover this manifold?



**Can we unfold this?**

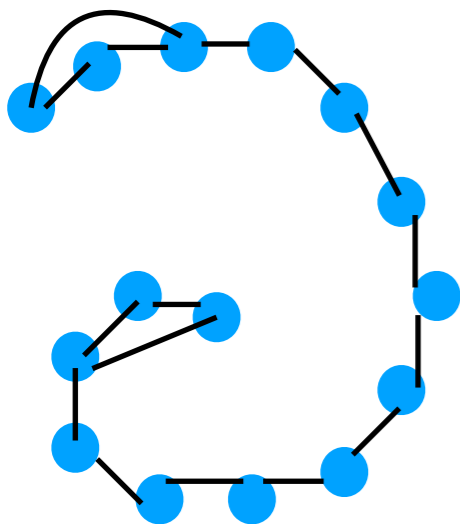
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors



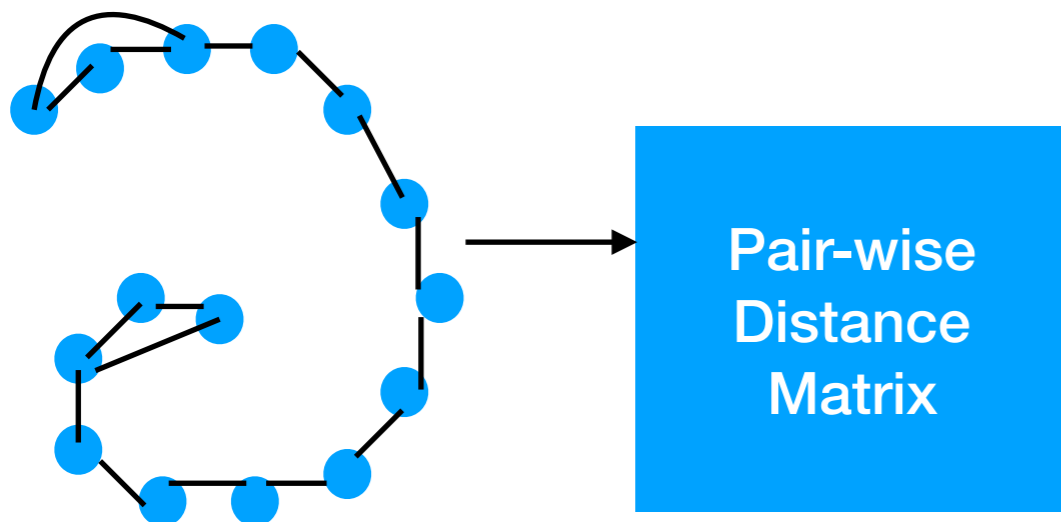
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph



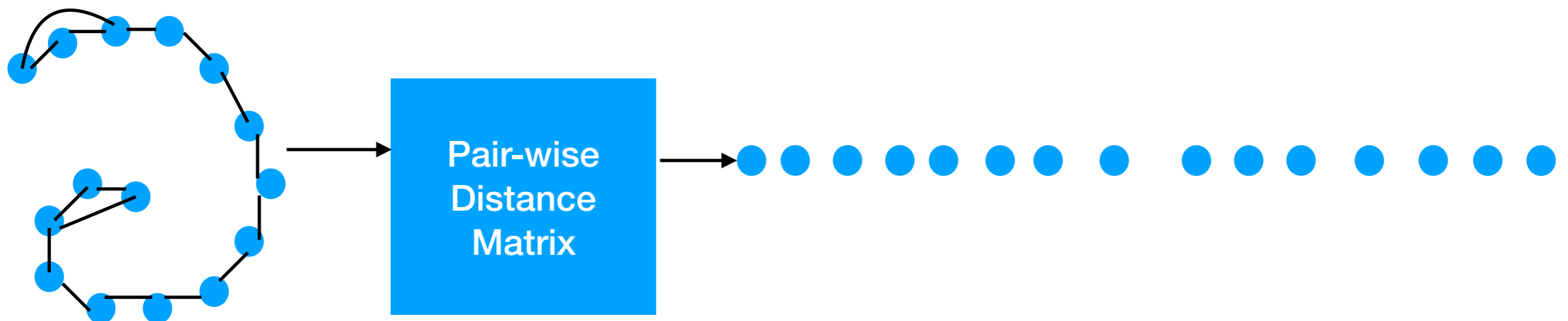
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph
- 3 For every pair of points  $A$  and  $B$ , distance between point  $A$  to  $B$  is shortest distance between  $A$  and  $B$  on graph



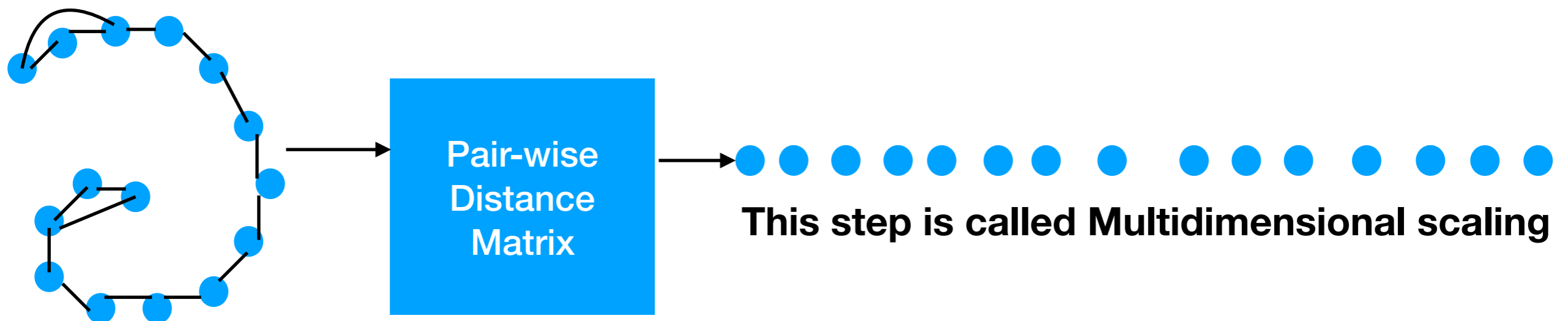
# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph
- 3 For every pair of points  $A$  and  $B$ , distance between point  $A$  to  $B$  is shortest distance between  $A$  and  $B$  on graph
- 4 Find points in low dimensional space such that distances between points in this space is equal to distance on graph.



# METHOD I: ISOMAP

- 1 For every point, find its ( $k$ -) Nearest Neighbors
- 2 Form the Nearest Neighbor graph
- 3 For every pair of points  $A$  and  $B$ , distance between point  $A$  to  $B$  is shortest distance between  $A$  and  $B$  on graph
- 4 Find points in low dimensional space such that distances between points in this space is equal to distance on graph.



# ISOMAP: PITFALLS

- ① If we don't take enough nearest neighbors, then graph may not be connected
- ② If we connect points too far away, points that should not be connected can get connected
- ③ There may not be a right number of nearest neighbors we should consider!



Demo